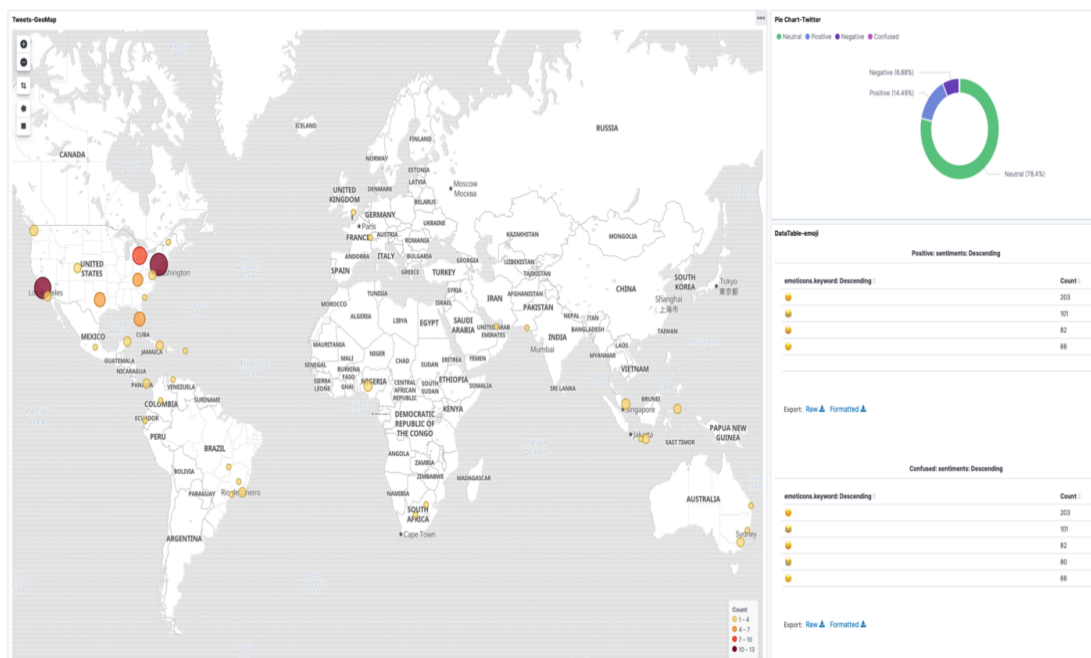# TwitterAPI with OpenSearch

> using knowledge of Cloud Computing Course which tough by professor Unubold Tumenbayar From Streaming Data to COVID-19 Twitter Analysis: Using AWS Lambda, Kinesis Firehose and Elasticsearch
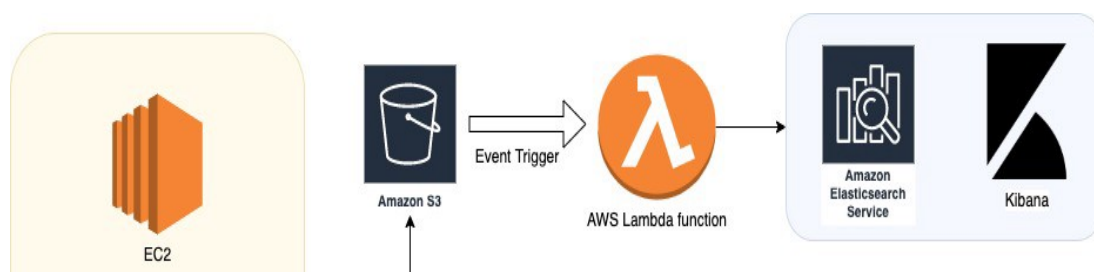
# From Streaming Data to COVID-19 Twitter Analysis: Using AWS Lambda, Kinesis Firehose and Elasticsearch
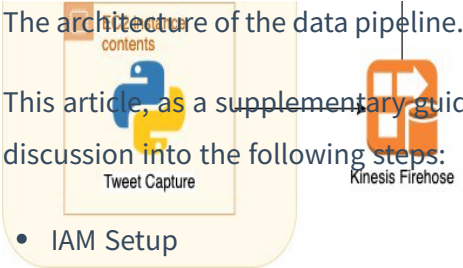
## Unleash the power of AWS and Elasticsearch



The dashboard of analyzing the real-time tweets.

The data pipeline follows the fashion of **near real-time streaming**, from digesting live Twitter data, to visualization. Most of the fragments to assemble the puzzle are from the AWS family: **AWS Kinesis Firehose**, **AWS S3 Bucket**, **AWS Lambda**, and **Amazon Elasticsearch Service**. The architecture is like this:

The architecture of the data pipeline.

This article, as a supplementary guideline to the *reference article*, will separate the discussion into the following steps:

- IAM Setup
- Creating the Amazon Elasticsearch Service cluster
- Configuring the AWS Kinesis Firehose and S3
- Creating the AWS Lambda function
- Code packaging and changes
- Kibana visualization and Twitter analysis

The code for **the Lambda function** and **Twitter capture program** has been uploaded in my public repo.

# IAM Setup

Before building the data pipeline, you need to have an AWS account and Twitter API keys and access tokens, which is also mentioned in the prerequisites of *the reference article*. Besides that, IAM roles are critical and must be set up correctly. Two roles are needed:

- Kinesis Firehose needs an IAM role with granted permissions to deliver stream data, which will be discussed in the section of Kinesis and S3 bucket.
- AWS Lambda needs permissions to access the S3 event trigger, add CloudWatch logs, and interact with Amazon Elasticserch Service.

Roles  >  lambda-s3-es-role

## Summary

Delete role

| | |
|---|---|
| **Role ARN** | arn:aws:iam::          :role/service-role/lambda-s3-es-role |
| **Role description** | Edit |
| **Instance Profile ARNs** | |
| **Path** | /service-role/ |
| **Creation time** | 2020-03-16 16:27 EDT |
| **Last activity** | 2020-04-16 13:11 EDT (Today) |
| **Maximum CLI/API session duration** | 1 hour Edit |

Permissions | Trust relationships | Tags | Access Advisor | Revoke sessions

▼ Permissions policies (3 policies applied)

Attach policies          ⊕ Add inline policy

| Policy name ▼ | Policy type ▼ | |
|---|---|---|
| ▶    AWSLambdaS3ExecutionRole-5227dcd9-d767-4527-aa70-0... | Managed policy | ✖ |

The IAM role, lambda-s3-es-role, for the Lambda function.

AWSLambdaBasicExecutionRole-9d8950ec-dfa7-4c4c-9621...

Shown as the above image, I attached three policies to the Lambda execution role *lambda-s3-es-role*. If you are not sure how to configure the policies, I attach these policies to the **repo** for reference.

AWSLambdaElasticsearchExecutionRole-3e5bdfc4-d778-4ff,...

▸ Permissions boundary (not set)

# Creating the Amazon Elasticsearch Service (ES) cluster

I assume that readers have followed through the steps in the *reference article* to create an Amazon ES domain at **Amazon ES home page**. For free-tier users, they can choose instance types as t2.micro or t2.small, and earn free 750 hours usage of Amazon ES. A few points need to note when creating the ES domain:

- No need to set up "dedicated master nodes".
- As a demo project, I choose public access in "network configuration".

Here is the thing. When configuring "Access policy", choosing "customer access policy", you need to add the following policy:

1. Select "ARN", and allow the lambda execution role *lambda-s3-es-role* to access ES service.

Modify the access policy for my-domain

Access policies control whether a request is accepted or rejected when it reaches the Amazon Elasticsearch Service domain. If you specify an account, user, or role in this policy, you must sign your requests. Learn more

Status   **Active**

Domain access policy   | Custom access policy           ▼ |

Allow or deny access by AWS account ID, account ARN, IAM user ARN, IAM role ARN, IPv4 address, or CIDR block.

| IAM ARN        ▼ |   arn:aws:iam::.............. |   Allow        ▼ |   Remove element |

Add element

Back   **Submit**

I leave a question here: **are the settings correct?** We will test it later.

# Configuring the AWS Kinesis Firehose and S3

Different from the *reference article*, I choose to create a Kinesis Firehose at the **Kinesis Firehose Stream console**. The steps are simple:

- Fill a name for the Firehose Stream
- Source: Direct PUT or other sources
- Destination: an S3 bucket, which is used to store data files (actually, tweets). Here you can choose an S3 bucket you have created or create a new one on the fly.
- Permissions.

As mentioned in the **IAM Section**, a Firehose Stream needs IAM roles to contain all necessary permissions. Click "Create new or choose", and choose to "create a new IAM Role", or use an existing one. The default policy would be attached and should meet the need.

Amazon Kinesis Firehose is requesting permission to use resources in your account

Click Allow to give Amazon Kinesis Firehose Read and Write access to resources in your account.

▼ Hide Details

Role Summary ❓

Role Description                    Create a new IAM Role          nd Resources
                                    firehose_delivery_role

IAM Role          ✓ Firehose_delivery_role_███████

Policy Name              Create a new Role Policy            ⬍

▶ View Policy Document

Kinesis Data Firehose Stream (KDF) and Kinesis Data Stream (KDS) may confuse people sometimes. KDF has extra features when delivering stream data. Source data is allowed to be transformed through a Lambda function during delivery to destination. My other **post** covers the usage of KDF.

# Creating the AWS Lambda function

AWS Lambda plays a central role in this pipeline. We will create a Lambda function to do the following jobs:

- Once a new data file created in the target S3 bucket, the Lambda function would be triggered.
- Data would be parsed with a designated structure, which agrees with the mapping for documents.
- The data would be loaded into the ES cluster.

To implement such a Lambda function at one stroke is hard. Dividing the cumbersome procedure into smaller steps, I first need to set up the Lambda environment correctly.

Create a function at **AWS Lambda home page**:

- Choosing Python 3.7 runtime.
- Choosing *lambda-s3-es-role* as the execution role.
- Keeping memory at 128 MB and set timeout as 2 min.
- Adding a trigger for S3. If any new file comes to the S3 bucket, the Lambda function would receive the event and get invocated.

Add trigger

Create an S3 trigger for Lambda function, note that prefix "debug" is used for debugging and can be substituted at your needs.

Now, we can test if the Lambda function could react to the S3 event. With sample code and configured Handler, we put a file into the S3 bucket *twitter-stream-sink*.



Create a Lambda function and put a test file into the S3 bucket. Note that the Handler name should match with the entry function in Function code.

On the "Monitoring" tab on the Lambda function panel, there is one dot appearing on the metrics graphs. Clicking "View logs in CloudWatch", we have the *CloudWatch* log for this invocation, and the log prints the source S3 bucket and the file name which we just put in.



The left window shows a dot in metrics graphs, denoting an invocation on the Lambda function. The right window shows the CloudWatch log details.

# Code packaging and changes

The *reference article* was published several years ago, so the project code needs updates.

The project code can be downloaded from the **repo**. The code directory contains 4 python files: *config.py, myhandle.py, tweet_utils.py, twitter_to_es.py*. To add necessary libraries into the project folder, just need to type the command:

```
pip install library_name -t .
```

The libraries needed to import to the project directory are such as:

```
requests
requests_aws4auth
elasticsearch
textblob=0.15
```

The Lambda function accepts a code package with zip format. Packaging the code directory with libraries needs this command:

```
zip -r ../your_package_name.zip * -x "*.git*"
```

The tutorial about the Lambda deployment package can be found in **AWS Lambda document**.

Now let us take a glimpse at each python file:

## myhandle.py

- Serve as the entry point for the Lambda function.
- Parse the event information, and obtain S3 file content in JSON format.

## twitter_to_es.py

- Add indices and mappings to the ES cluster.
- Adopt the **bulk** way to load parsed data into the ES cluster.
- Authorize the requests sent to the ES cluster.

## tweet_utils.py

- Act as a helper module.
- Parse tweets into the structured dictionary.
- Analyze the sentiments over tweets using TextBlob.

# config.py

- Act as the shared configuration.

Compared with the **original code** in the *reference article*, I made some code changes:

1. Add extra libraries into the package, requests_aws4auth, requests.

2. Port source code from Python 2 to Python 3.

3. Fix the bugs because Elasticsearch and its Python client library have the incompatible issues with previous versions.

To port source code from Python 2 to 3, we can use the library `2to3`.

```
[hzhong@hzhongs-MacBook-Pro twitter-s3-to-es % 2to3 tweet_utils.py
RefactoringTool: Skipping optional fixer: buffer
RefactoringTool: Skipping optional fixer: idioms
RefactoringTool: Skipping optional fixer: set_literal
RefactoringTool: Skipping optional fixer: ws_comma
RefactoringTool: Refactored tweet_utils.py
--- tweet_utils.py        (original)
+++ tweet_utils.py        (refactored)
@@ -101,7 +101,7 @@


 def _sentiment_analysis_by_emoticons(tweet):
-        for sentiment, emoticons_icons in emoticons.iteritems():
+        for sentiment, emoticons_icons in emoticons.items():
            matched_emoticons = re.findall(emoticons_icons, tweet['text'].encode('utf-8'))
            if len(matched_emoticons) > 0:
                tweet['emoticons'].extend(matched_emoticons)
@@ -128,7 +128,7 @@
 def get_tweet(doc):
        tweet = {}
        tweet[id_field] = doc[id_field]
-        tweet['hashtags'] = map(lambda x: x['text'],doc['entities']['hashtags'])
+        tweet['hashtags'] = [x['text'] for x in doc['entities']['hashtags']]
        tweet['coordinates'] = doc['coordinates']
        tweet['timestamp_ms'] = doc['timestamp_ms']
        tweet['text'] = doc['text']
RefactoringTool: Files that need to be modified:
RefactoringTool: tweet_utils.py
```

An example of using *2to3 command to port Python 2 code to Python 3.*

The incompatible issues have been fixed:

- Since the release of Elasticsearch 7.0.0, **mapping types are removed**.
- Elasticsearch's **Python Client** also experiences changes from previous releases, particularly the usage of the *bulk* method.

To get familiar with the Elasticsearch's Python client, you can open a *Jupiter notebook* to test the connection with the ES cluster.

After looking into the code, we need to package the code. To test the Lambda function, we put a captured twitter file into the S3 bucket, and see if the tweets are parsed correctly and loaded into the ES cluster.



The left window shows the panel of the Lambda function, and the right window shows putting a sample data file with tweets into S3.

If the data is loaded into the ES cluster successfully, we can use Kibana's "Discover" function to check it.



The left window shows the trace logs during invocating the Lambda function, and the right window visualizes the tweets loaded into ES.

Apart from the code changes on the Lambda function, I use a Python program running on an AWS EC2 instance to capture tweets, which can be found **here**.

The *reference article* contains a **node.js program** to capture tweets, and either of them does the same job.



Running the Python tweet capture program

One thing to note here. When I implemented the pipeline and tried to test to load data into the ES cluster, I met authentication errors on Lambda invocation and Kibana:



When I put a data file into the S3 bucket, the Lambda function reported such an error.



When I tried to access Kibana after setting up the ES cluster.

To find out the reason, we need to go to the page specifying the **Identity and Access Management** of Amazon ES.

> The primary appeal of IP-based policies is that they allow unsigned requests to an Amazon ES domain, which lets you use clients like **curl** and **Kibana** or access the domain through a proxy server.

*All requests to the Amazon ES configuration API must be signed*. In the above errors, even though the Lambda execution role was added, the requests were unsigned and rejected, especially HTTP requests. To resolve this problem, we can add IP-based policies, or adding signs using **AWS SDK or requests**. So we need to add an IP-based policy into the ES access policies:



# Kibana visualization and Twitter analysis

Once we start running the twitter capture program, loads of tweets would be flown into S3 and the Lambda function would handle the data files. The best way to reflect the data is through the visualization utility *Kibana*, provided in Elasticsearch.

All data loaded into Elasticsearch need to be assigned with indices, and thus Kibana can use *index patterns* to retrieve the data. In *twitter_to_es.py*, tweets are indexed with "twitter". Now we can create an index pattern "twitter*" and start discovering the data in Kibana.

Create the index pattern "twitter*" to match all indices starting with "twitter".



The "timestamp_ms" field is specified when we add the mapping in tweet_utils.py.

After the creation of the index pattern, we can explore the data by choosing the "Discover" button on the left-sidebar and the data can be presented in time series.

In this article, I pick the hashtags "**#StayHome**" and "**#SocialDistancing**" to mine Twitter. Like the *reference article*, I create a dashboard to visualize the tweets. There are three visualizations in the dashboard:

- A coordinate map to demonstrate the geographical distribution of the tweets, only valid if tweets contain location information.
- A pie chart to demonstrate sentiment popularity when users send tweets, including three types of emotions, positive, neutral and negative.
- A data table to count the emojis contained in each sentiment, and only the top-5 emojis are listed.

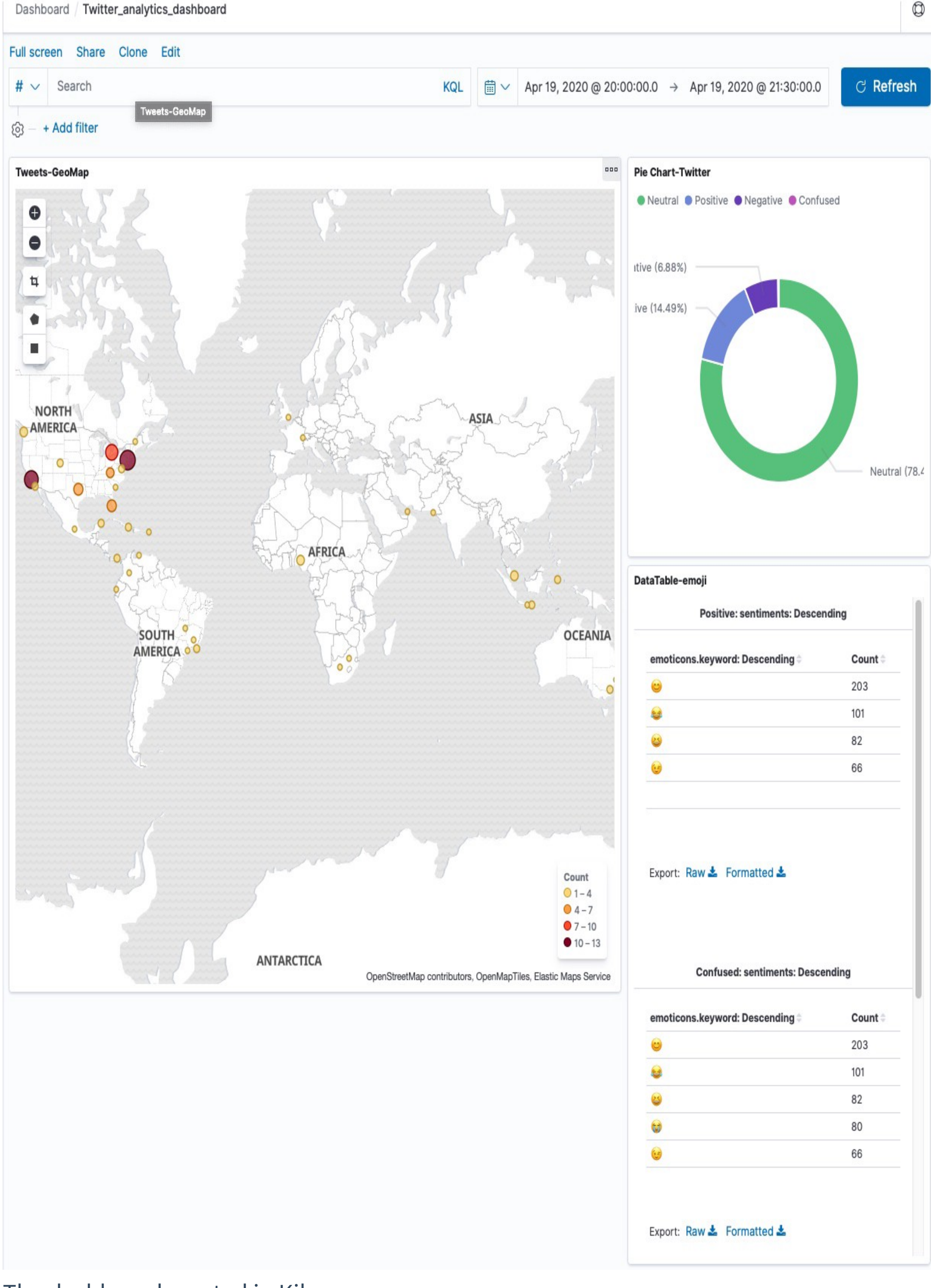


The dashboard created in Kibana.

# Conclusion

AWS Lambda and Elasticsearch are pretty powerful technologies and this post may just demonstrate one case among the application scenarios. Besides real-time data processing, Lambda can be integrated with ETL (Extract, Transform, Load) and application backends. Elasticsearch has established a reputation on logging/log analysis and full-text searching.

I hope you can find fun when done the reading and being able to fiddle with the big data technologies. This is how big data fascinates me.

# Conclusion